

## CLAIMS

1. A method comprising:  
separating at least a portion of an audio signal into a plurality of frames;  
extracting line spectrum pairs from each of the plurality of frames; and  
using at least the line spectrum pairs to classify at least the portion as either  
speech or non-speech.

2. A method as recited in claim 1, wherein the using comprises:  
generating an input Gaussian Model corresponding to the plurality of  
frames based on the extracted line spectrum pairs;  
comparing the input Gaussian Model to a Vector Quantization codebook  
including a plurality of trained Gaussian Models;  
identifying one of the plurality of trained Gaussian Models that is closest to  
the input Gaussian Model;  
determining a distance between the input Gaussian Model and the closest  
trained Gaussian Model; and  
classifying at least the portion as speech if the distance is less than a  
threshold value.

3. A method as recited in claim 1, wherein the using comprises:  
generating an input Gaussian Model corresponding to the plurality of  
frames based on the extracted line spectrum pairs;  
identifying one of the plurality of trained Gaussian Models that is closest to  
the input Gaussian Model;

1 determining a distance between the input Gaussian Model and the closest  
2 trained Gaussian Model; and

3 classifying at least the portion as non-speech if the distance is greater than a  
4 first threshold value.

5  
6 4. A method as recited in claim 3, further comprising:

7 determining an energy distribution of the plurality of frames in a first  
8 bandwidth; and

9 classifying at least the portion as non-speech if the distance is greater than a  
10 second threshold value and the energy distribution of the plurality of frames in the  
11 first bandwidth is less than a third threshold value, wherein the second threshold  
12 value is less than the first threshold value.

13  
14 5. A method as recited in claim 4, further comprising:

15 determining an energy distribution of the plurality of frames in a second  
16 bandwidth; and

17 classifying at least the portion as speech if the distance is less than the  
18 second threshold value and the energy distribution of the plurality of frames in the  
19 second bandwidth is greater than a fourth threshold value.

20  
21 6. A method as recited in claim 5, further comprising otherwise  
22 classifying at least the portion as speech.

1           7.    A method as recited in claim 2, further comprising:  
2           extracting a high zero crossing rate ratio feature from the plurality of  
3 frames;  
4           extracting a low short time energy ratio feature from the plurality of frames;  
5           extracting a spectrum flux feature from the plurality of frames;  
6           pre-classifying the portion as speech or non-speech based at least in part on  
7 an average zero crossing rate, the high zero crossing rate ratio, the low short time  
8 energy ratio, and the spectrum flux features;  
9           using a first value as the threshold value if the portion is pre-classified as  
10 speech; and  
11           using a second value as the threshold value if the portion is pre-classified as  
12 non-speech, wherein the second value is greater than the first value.

13  
14           8.    One or more computer-readable memories containing a computer  
15 program that is executable by a processor to perform the method recited in claim  
16 1.

17  
18        ~~9.~~    A method comprising:  
19           separating at least a portion of an audio signal into a plurality of frames;  
20           extracting a periodicity feature from the plurality of frames; and  
21           using at least the periodicity feature to classify at least the portion as either  
22 music or environment sound.  
23  
24  
25

1           **10.**    A method as recited in claim 9, wherein the periodicity feature  
2 comprises a noise frame ratio that identifies a ratio of noise frames to non-noise  
3 frames in the plurality of frames.

4  
5           **11.**    A method as recited in claim 10, further comprising classifying at  
6 least the portion as environment sound if the noise frame ratio exceeds a threshold  
7 value.

8  
9           **12.**    A method as recited in claim 10, further comprising:  
10           extracting, from the plurality of frames, a band periodicity for each of a  
11 plurality of bands of the audio signal and a full band periodicity that is a  
12 concatenation of the band periodicities for each of the plurality of bands; and  
13           classifying at least the portion as environment sound if the first band  
14 periodicity is less than a first threshold or the second band is less than a second  
15 threshold.

16  
17           **13.**    A method as recited in claim 9, wherein the periodicity feature  
18 comprises a band periodicity for each of a plurality of bands of the audio signal.

19  
20           **14.**    A method as recited in claim 13, further comprising:  
21           extracting a full band periodicity from the plurality of frames that is a  
22 concatenation of the band periodicities for each of the plurality of bands; and  
23           classifying at least the portion as environment sound if the full band  
24 periodicity exceeds a threshold value.  
25

1           15. A method as recited in claim 9, further comprising extracting a  
2 spectrum flux feature from the plurality of frames, and wherein the using  
3 comprises using at least the periodicity feature and the spectrum flux feature to  
4 classify at least the portion as either music or environment sound.

5  
6           16. A method as recited in claim 15, wherein the spectrum flux feature  
7 is extracted by determining a Fast Fourier Transform for each of the plurality of  
8 frames and calculating a difference in the Fast Fourier Transforms for successive  
9 frames.

10  
11           17. A method as recited in claim 15, further comprising:  
12 extracting, from the plurality of frames, an energy distribution in a band of  
13 the audio signal; and  
14 classifying at least the portion as environment sound if the band energy  
15 distribution is less than a first threshold or the spectrum flux exceeds a second  
16 threshold.

17  
18           18. A method as recited in claim 9, wherein the periodicity feature  
19 comprises a band periodicity for each of a plurality of bands of the audio signal,  
20 and further comprising:

21 extracting, from the plurality of frames, a spectrum flux feature;  
22 extracting, from the plurality of frames, an energy feature indicating an  
23 amount of energy in at least one band of the portion; and  
24  
25

1 classifying at least the portion as environment sound if the amount of  
2 energy is less than a first threshold and the spectrum flux is less than a second  
3 threshold.

4  
5 19. One or more computer-readable memories containing a computer  
6 program that is executable by a processor to perform the method recited in claim  
7 9.

8  
9 ~~20.~~ A method comprising:  
10 separating at least a portion of an audio signal into a plurality of frames;  
11 extracting a periodicity feature for each of the plurality of frames; and  
12 using at least the periodicity feature to classify the plurality of frames as  
13 either music with vocals or music without vocals.

14  
15 21. A method as recited in claim 20, wherein the periodicity feature  
16 comprises a band periodicity for each of a plurality of bands of the audio signal.

17  
18 22. A method as recited in claim 21, further comprising classifying at  
19 least the portion as music with vocals if the band periodicity of at least one of the  
20 plurality of bands is greater than a first threshold and less than a second threshold.  
21  
22  
23  
24  
25

1           23.     A method as recited in claim 22, further comprising classifying at  
2     least the portion as environment sound if the band periodicity of each of the  
3     plurality of bands is less than the second threshold, and otherwise classifying at  
4     least the portion as music without vocals.

5  
6           24.     One or more computer-readable memories containing a computer  
7     program that is executable by a processor to perform the method recited in claim  
8     20.

9  
10          ~~25.~~    A method for determining when a speaker changes, the method  
11     comprising:

12               separating at least a portion of an audio signal into a plurality of frames;  
13               extracting line spectrum pairs from each of the plurality of frames; and  
14               determining when a speaker of the audio signal changes based at least in  
15     part on the line spectrum pairs.

16  
17          26.     A method as recited in claim 25, wherein the determining  
18     comprises:

19               calculating a difference between line spectrum pairs for successive frames  
20     of the plurality of frames;

21               if the difference between two line spectrum pairs exceeds a threshold value,  
22     then determining that the speaker has changed, otherwise determining that the  
23     speaker has not changed.

1           **27.**   One or more computer-readable memories containing a computer  
2 program that is executable by a processor to perform the method recited in claim  
3 25.

4  
5           **28.**   An apparatus comprising:  
6           a line spectrum pair (LSP) analyzer to extract line spectrum pairs from a  
7 portion of an audio signal; and  
8           a speech discriminator, communicatively coupled to the LSP analyzer, to  
9 classify the portion of the audio signal as either speech or non-speech based at  
10 least in part on the LSP analyzer.

11  
12           **29.**   An apparatus as recited in claim 28, further comprising:  
13           a distance calculator, communicatively coupled to the LSP analyzer, to  
14 determine a distance between at least one of the trained Gaussian Models and an  
15 input Gaussian Model based on the extracted line spectrum pairs; and

16           wherein the speech discriminator is further to classify the portion of the  
17 audio signal as either speech or non-speech based at least in part on the distance  
18 between the at least one of the trained Gaussian Models and the input Gaussian  
19 Model.

20  
21           **30.**   An apparatus as recited in claim 28, further comprising:  
22           a Fast Fourier Transform (FFT) analyzer to extract Fast Fourier Transform  
23 features from the portion of the audio signal;  
24  
25



1 an energy distribution calculator, communicatively coupled to both the FFT  
2 analyzer and the speech discriminator, to determine an energy distribution of the  
3 portion of the audio signal in at least one bandwidth; and

4 wherein the speech discriminator is further to classify the portion of the  
5 audio signal as either speech or non-speech based at least in part on the energy  
6 distribution of the portion of the audio signal in the at least one bandwidth.

7  
8 **31.** An apparatus comprising:

9 a band periodicity calculator to determine a periodicity of each of a  
10 plurality of bands of a portion of an audio signal; and

11 a discriminator, communicatively coupled to the band periodicity  
12 calculator, to classify the portion of the audio signal as music or environment  
13 sound based at least in part on the periodicity of one of the plurality of bands.

14  
15 **32.** An apparatus as recited in claim 31, further comprising:

16 a noise frame ratio calculator, communicatively coupled to the  
17 discriminator, to determine a noise frame ratio of the portion of the audio signal;  
18 and

19 wherein the discriminator is to classify the portion of the audio signal as  
20 music or environment sound based at least in part on the periodicity of one of the  
21 plurality of bands and on the noise frame ratio of the portion.

22  
23 **33.** An apparatus as recited in claim 31, further comprising:

24 a spectrum flux analyzer, communicatively coupled to the discriminator, to  
25 determine a spectrum flux of the portion of the audio signal; and

1 wherein the discriminator is to classify the portion of the audio signal as  
2 music or environment sound based at least in part on the periodicity of one of the  
3 plurality of bands and on the spectrum flux of the portion.

4  
5 **34.** A method comprising:  
6 receiving an audio signal;  
7 separating the audio signal into a plurality of portions; and  
8 classifying each of the plurality of portions, based at least in part on  
9 periodicity features of the portion, as one of: speech, music, silence, and  
10 environment sound.

11  
12 **35.** A method as recited in claim 34, wherein the periodicity features  
13 include a noise frame ratio that identifies a ratio of noise frames to non-noise  
14 frames in the plurality of frames.

15  
16 **36.** A method as recited in claim 35, wherein the classifying comprises  
17 classifying at least the portion as environment sound if the noise frame ratio  
18 exceeds a threshold value.

19  
20 **37.** A method as recited in claim 34, further comprising:  
21 extracting, from the plurality of frames, a band periodicity for each of a  
22 plurality of bands of the audio signal and a full band periodicity that is a  
23 concatenation of the band periodicities for each of the plurality of bands; and  
24 wherein the classifying comprises classifying at least the portion as  
25 environment sound if a band periodicity of a first of the plurality of bands is less

1 than the first threshold a band periodicity of a second of the plurality of bands is  
2 less than the second threshold.

3  
4 **38.** A method as recited in claim 34, wherein the periodicity features  
5 include a band periodicity for each of a plurality of bands of the audio signal.

6  
7 **39.** A method as recited in claim 38, further comprising:  
8 extracting a full band periodicity from the plurality of frames that is a  
9 concatenation of the band periodicities for each of the plurality of bands; and  
10 wherein the classifying comprises classifying at least the portion as  
11 environment sound if the full band periodicity exceeds a threshold value.

12  
13 **40.** A method as recited in claim 34, further comprising:  
14 extracting a spectrum flux feature from the plurality of frames; and  
15 wherein the classifying comprises classifying at least the portion as either  
16 music or environment sound based at least in part on the periodicity feature and  
17 the spectrum flux feature.

18  
19 **41.** A method as recited in claim 34, further comprising:  
20 extracting line spectrum pairs from each of the plurality of frames;  
21 generating an input Gaussian Model corresponding to the plurality of  
22 frames based on the extracted line spectrum pairs;  
23 comparing the input Gaussian Model to a Vector Quantization codebook  
24 including a plurality of trained Gaussian Models;

1 identifying one of the plurality of trained Gaussian Models that is closest to  
2 the input Gaussian Model;

3 determining a distance between the input Gaussian Model and the closest  
4 trained Gaussian Model; and

5 classifying at least the portion as speech if the distance is less than a  
6 threshold value.

7  
8 **42.** A method as recited in claim 34, further comprising:

9 extracting line spectrum pairs from each of the plurality of frames;

10 generating an input Gaussian Model corresponding to the plurality of  
11 frames based on the extracted line spectrum pairs;

12 identifying one of the plurality of trained Gaussian Models that is closest to  
13 the input Gaussian Model;

14 determining a distance between the input Gaussian Model and the closest  
15 trained Gaussian Model; and

16 classifying at least the portion as one of music, silence, or environment  
17 sound if the distance is greater than a first threshold value.

18  
19 **43.** A method as recited in claim 42, further comprising:

20 determining an energy distribution of the plurality of frames in a first  
21 bandwidth; and

22 classifying at least the portion as one of music, silence, or environment  
23 sound if the distance is greater than a second threshold value and the energy  
24 distribution of the plurality of frames in the first bandwidth is less than a third  
25

1 threshold value, wherein the second threshold value is less than the first threshold  
2 value.

3 . 44. A method as recited in claim 43, further comprising:  
4 determining an energy distribution of the plurality of frames in a second  
5 bandwidth; and  
6 classifying at least the portion as one of music, silence, or environment  
7 sound if the distance is greater than a fourth threshold value and the energy  
8 distribution of the plurality of frames in the second bandwidth is less than a fifth  
9 threshold value, wherein the fourth threshold value is less than the first threshold  
10 value.

11  
12 45. A method as recited in claim 44, further comprising otherwise  
13 classifying at least the portion as speech.

14  
15 46. One or more computer-readable memories containing a computer  
16 program that is executable by a processor to perform the method recited in claim  
17 34.

18  
19 47. One or more computer-readable media having stored thereon a  
20 computer program to classify a portion of an audio signal as speech, music,  
21 silence, or environment sound, wherein the computer program, when executed by  
22 one or more processors, causes the one or more processors to perform acts  
23 including:

24 (a) analyzing line spectrum pair features of the portion to determine if  
25 the portion is speech;

1 (b) analyzing energy features of the portion to determine if the portion is  
2 silence;

3 (c) analyzing periodicity features of the portion to determine if the  
4 portion is music or environment sound; and

5 (d) classifying the portion as speech, music, silence, or environment  
6 sound based on at least one of the analyzing acts (a)-(c).

7  
8 48. One or more computer-readable media as recited in claim 47,  
9 wherein the computer program is further to cause the one or more processors to  
10 perform the acts (a) – (d) in the order (a), then (b), then (c), then (d).

11  
12 49. One or more computer-readable media as recited in claim 48,  
13 wherein the computer program is further to cause the one or more processors to  
14 perform act (b) only if act (a) results in a determination that the portion is not  
15 speech.

16  
17 50. One or more computer-readable media as recited in claim 48,  
18 wherein the computer program is further to cause the one or more processors to  
19 perform act (c) only if act (b) results in a determination that the portion is not  
20 silence.